

Я знаю, что ты знаешь, что я знаю: интерактивная рациональность и байесовская прагматика *

Виталий Долгоруков

1 Он знает, что я знаю, что он знает

«Он сказал, что p ; нет оснований считать, что он не соблюдает постулаты или по крайней мере принцип кооперации; он не мог сказать p , если бы он не считал, что q ; он знает (и знает, что я знаю, что он знает), что я могу понять необходимость предположения о том, что он думает, что q ; он хочет, чтобы я думал – или хотя бы готов позволить мне думать – что q : итак, он имплицировал, что q » [Грайс 1985, 227-228]

2 Модель Iterated Best Responce (IBR)

1) Franke 2009; Jäger 2013

2) Скалярные импликатуры

Я съел несколько конфет \leadsto Я съел не все конфеты

2.1 Сигнальная игра

$G = \langle \{S, R\}, W, Pr, F, A, \mathcal{S}, \mathcal{R}, \rho, \sigma, U_S, U_R \rangle$, где

- $\{S, R\}$ - Говорящий и Слушающий
- $W = \{w_{\forall}, w_{\exists \rightarrow \forall}\}$ - множество состояний (возможных миров)
- $A = \{a_{\forall}, a_{\exists \rightarrow \forall}\}^1$ - множество действий (интерпретаций)
- $Pr \in \Delta(W)^2$ - распределение вероятностей над W
- $F = \{\varphi_{some}, \varphi_{all}\}$ - множество сообщений (форм)
- $\mathcal{S} = F^W$ - множество стратегий Говорящего
- $\mathcal{R} = A^F$ - множество стратегий Слушающего
- $U_S = U_R : W \times A \mapsto \mathbb{R}$ - платежная функция
- $\rho \in \Delta(A)^F$ - представление о поведении Слушающего

* «Динамический поворот в логической семантике» (Научный фонд НИУ ВШЭ, 15-05-0005, НУГ «Формальная философия»)

¹В общем виде: $|A| = |W|$ и $\forall w_i \in W \exists a_j \in A i = j$

² $\Delta(X) = \{f : X \mapsto [0, 1] \mid \sum_{x \in X} f(x) = 1\}$.

- $\sigma \in \Delta(F)^W$ - представление о поведении Говорящего
- $\|\cdot\|_W : F \mapsto \mathcal{P}(W)$, $\|\cdot\|_A : F \mapsto \mathcal{P}(A)$ - функции интерпретации

2.2 Семантическая компетентность

2.2.1 Принцип кооперации

Будем считать, что Говорящий и Слушающий соблюдают Принцип Кооперации.

$$U_S(w_i, a_j) = U_R(w_i, a_j) = \begin{cases} 1, & \text{if } i = j \\ -1, & \text{else} \end{cases} \quad (1)$$

| $Pr(w)$ | w | $U(w, a_{\forall})$ | $U(w, a_{\exists \rightarrow \forall})$ | $\ \varphi_{some}\ _W$ | $\ \varphi_{all}\ _W$ |
|---------|-----------------------------------|---------------------|---|------------------------|-----------------------|
| $1/2$ | w_{\forall} | (1; 1) | (-1; -1) | ✓ | ✓ |
| $1/2$ | $w_{\exists \rightarrow \forall}$ | (-1; -1) | (1; 1) | ✓ | × |

Таблица 1: Скалярные импликатуры

2.3 Оптимальные стратегии

2.3.1 Оптимальная стратегия для Говорящего

Ожидаемая полезность для Говорящего, который находясь в мире w , посылает сообщение φ (с учетом представления ρ):

$$EU_S(w, \varphi, \rho) = \sum_{a' \in A} \rho(\varphi, a') \times U_S(w, a') \quad (2)$$

Наилучший ответ для Говорящего в мире w (с учетом представления ρ), (то есть, то сообщение, отправка которого в мире w максимизирует ожидаемую

полезность):

$$br_S(w, \rho) = \{\varphi \mid EU_S(w, \varphi, \rho) = \max_{\varphi' \in F} EU_S(w, \varphi', \rho)\} \quad (3)$$

Наилучший ответ Говорящего на его представление ρ :

$$BR_S(\rho) = \{s \in \mathcal{S} \mid \forall w : s(w) \in br_S(w, \rho)\} \quad (4)$$

То есть, та стратегия Говорящего, которая для любого возможного мира максимизирует ожидаемую полезность.

| w | φ | $U_S(w, a) / \rho(\varphi, a)$ | a_{\forall} | $a_{\exists \neg \forall}$ | $EU_S(w, \varphi, \rho)$ |
|----------------------------|------------------|------------------------------------|---------------|----------------------------|--------------------------|
| w_{\forall} | φ_{some} | $U_S(w_{\forall}, a)$ | 1 | -1 | 0 |
| | | $\rho(\varphi_{some}, a)$ | 1/2 | 1/2 | |
| | φ_{all} | $U_S(w_{\forall}, a)$ | 1 | -1 | 1 |
| | | $\rho(\varphi_{all}, a)$ | 1 | 0 | |
| $w_{\exists \neg \forall}$ | φ_{some} | $U_S(w_{\exists \neg \forall}, a)$ | -1 | 1 | 0 |
| | | $\rho(\varphi_{some}, a)$ | 1/2 | 1/2 | |
| | φ_{all} | $U_S(w_{\exists \neg \forall}, a)$ | -1 | 1 | -1 |
| | | $\rho(\varphi_{all}, a)$ | 1 | 0 | |

Таблица 2: $EU_S(w, \varphi, \rho)$

2.3.2 Оптимальная стратегия для Слушающего

$$EU_R(\varphi, a, \sigma) = \sum_{w' \in W} \mu(w' | \varphi) \times U_R(w', a) \quad (5)$$

$$br_R(\varphi, \sigma) = \{a \mid EU_R(\varphi, a, \sigma) = \max_{a' \in A} EU_R(\varphi, a', \sigma)\} \quad (6)$$

$$BR_R(\sigma) = \{r \in R \mid \forall \varphi : r(\varphi) \in br_R(\varphi, \sigma)\} \quad (7)$$

2.4 IBR-последовательность

$$\mathcal{S}_0 = \{s \in \mathcal{S} \mid \forall w \in W : w \in \|s(w)\|_W\} \quad (8)$$

$$\mathcal{R}_0 = \{r \in \mathcal{R} \mid \forall \varphi \in F : r(\varphi) \in \|\varphi\|_A\} \quad (9)$$

$$\sigma_{k+1}(w, \varphi) = \frac{|\{s \in (S)_k \mid s(w) = \varphi\}|}{|(S)_k|} \quad (10)$$

$\mu(w|\varphi)$ представление Слушающего, получившего сигнал φ , о том, что этот сигнал был отправлен из состояния w . Это представления получается по формуле Байесовской вероятности из Pr и σ_k .

$$\mu_k(w|\varphi) = \frac{\sigma_k(\varphi, w) \times Pr(w)}{\sum_{w' \in W} \sigma_k(\varphi, w') \times Pr(w')} \quad (11)$$

$$\rho_{k+1}(\varphi, a) = \frac{|\{r \in (R)_k \mid r(\varphi) = a\}|}{|(R)_k|} \quad (12)$$

$$\mathcal{R}_k = BR_R(\mu_k) \quad (13)$$

$$\mathcal{S}_k = BR_S(\rho_k) \quad (14)$$

2.4.1 $\mathcal{R}_0 \mapsto \rho_1 \mapsto \mathcal{S}_1 \mapsto \sigma_2/\mu_2 \mapsto \mathcal{R}_2 \dots$

$$\mathcal{R}_0 = \left\{ \left[\begin{array}{l} \varphi_{all} \mapsto a_{\forall} \\ \varphi_{some} \mapsto a_{\forall} \end{array} \right], \left[\begin{array}{l} \varphi_{all} \mapsto a_{\forall} \\ \varphi_{some} \mapsto a_{\exists \neg \forall} \end{array} \right] \right\} \quad (15)$$

$$\begin{array}{c|cc} \rho_1(\varphi, a) & a_{\forall} & a_{\exists \neg \forall} \\ \varphi_{some} & 1/2 & 1/2 \\ \varphi_{all} & 0 & 1 \end{array} \quad (16)$$

$$\mathcal{S}_1 = BR_S(\rho_1) = \left\{ \left[\begin{array}{l} w_{\forall} \mapsto \varphi_{all} \\ w_{\exists \neg \forall} \mapsto \varphi_{some} \end{array} \right] \right\} \quad (17)$$

$$\begin{array}{c|cc|cc} \sigma_2(w, \varphi) & \varphi_{all} & \varphi_{some} & \mu_2(w|\varphi) & \varphi_{all} & \varphi_{some} \\ w_{\forall} & 1 & 0 & w_{\forall} & 1 & 0 \\ w_{\exists \neg \forall} & 0 & 1 & w_{\exists \neg \forall} & 0 & 1 \end{array} \quad (18)$$

$$\mathcal{R}_2 = BR_R(\mu_2) = \left\{ \left[\begin{array}{l} \varphi_{all} \mapsto a_{\forall} \\ \varphi_{some} \mapsto a_{\exists \neg \forall} \end{array} \right] \right\} \quad (19)$$

$$\begin{array}{c|cc} \rho_3(\varphi, a) & a_{\forall} & a_{\exists \neg \forall} \\ \varphi_{some} & 0 & 1 \\ \varphi_{all} & 1 & 0 \end{array} \quad (20)$$

$$\mathcal{S}_3 = BR_S(\rho_3) = \left\{ \left[\begin{array}{l} w_{\forall} \mapsto \varphi_{all} \\ w_{\exists \neg \forall} \mapsto \varphi_{some} \end{array} \right] \right\} \quad (21)$$

$\mathcal{S}_3 = \mathcal{S}_1$ – последовательность останавливается

2.4.2 $S_0 \mapsto \sigma_1/\mu_1 \mapsto \mathcal{R}_1 \mapsto \rho_2 \mapsto S_2 \dots$

$$\mathcal{S}_0 = \left\{ \left[\begin{array}{l} w_{\forall} \mapsto \varphi_{all} \\ w_{\exists-\forall} \mapsto \varphi_{some} \end{array} \right], \left[\begin{array}{l} w_{\forall} \mapsto \varphi_{some} \\ w_{\exists-\forall} \mapsto \varphi_{some} \end{array} \right] \right\} \quad (22)$$

$$\begin{array}{c|cc} \sigma_1(w, \varphi) & \varphi_{all} & \varphi_{some} \\ \hline w_{\forall} & 1/2 & 1/2 \\ w_{\exists-\forall} & 0 & 1 \end{array} \quad \begin{array}{c|cc} \mu_1(w|\varphi) & \varphi_{all} & \varphi_{some} \\ \hline w_{\forall} & 1 & 1/3 \\ w_{\exists-\forall} & 0 & 2/3 \end{array} \quad (23)$$

$$\mathcal{R}_1 = BR_R(\mu_0) = \left\{ \left[\begin{array}{l} \varphi_{all} \mapsto a_{\forall} \\ \varphi_{some} \mapsto a_{\exists-\forall} \end{array} \right] \right\} \quad (24)$$

$$\begin{array}{c|cc} \rho_2(\varphi, a) & a_{\forall} & a_{\exists-\forall} \\ \hline \varphi_{some} & 0 & 1 \\ \varphi_{all} & 1 & 0 \end{array} \quad (25)$$

$$\mathcal{S}_2 = BR_S(\rho_2) = \left\{ \left[\begin{array}{l} w_{\forall} \mapsto \varphi_{all} \\ w_{\exists-\forall} \mapsto \varphi_{some} \end{array} \right] \right\} \quad (26)$$

$$\begin{array}{c|cc} \sigma_3(w, \varphi) & \varphi_{all} & \varphi_{some} \\ \hline w_{\forall} & 1 & 0 \\ w_{\exists-\forall} & 0 & 1 \end{array} \quad \begin{array}{c|cc} \mu_3(w|\varphi) & \varphi_{all} & \varphi_{some} \\ \hline w_{\forall} & 1 & 0 \\ w_{\exists-\forall} & 0 & 1 \end{array} \quad (27)$$

$$\mathcal{R}_3 = BR_R(\mu_2) = \left\{ \left[\begin{array}{l} \varphi_{all} \mapsto a_{\forall} \\ \varphi_{some} \mapsto a_{\exists-\forall} \end{array} \right] \right\} \quad (28)$$

$\mathcal{R}_3 = \mathcal{R}_1$ – последовательность останавливается.

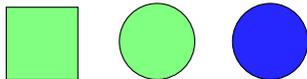
Таким образом, и Слушающий, и Говорящий приходят к выводу, что оптимальным вариантом поведения будет следующий профиль стратегий:

$$\left\langle \left[\begin{array}{l} w_{\forall} \mapsto \varphi_{all} \\ w_{\exists-\forall} \mapsto \varphi_{some} \end{array} \right], \left[\begin{array}{l} \varphi_{all} \mapsto a_{\forall} \\ \varphi_{some} \mapsto a_{\exists-\forall} \end{array} \right] \right\rangle \quad (29)$$

3 Модель Rational Speech Act (RSA)

Frank, Goodman 2012; Qing, Franke 2015

3.1 Эксперимент



Порождение: выбрать выражение (square, green, circle, blue)

Интерпретация: выбрать фигуру

| Referent | Property | | | |
|----------|----------|-------|--------|------|
| | square | green | circle | blue |
| | 1 | 0.5 | 0 | 0 |
| | 0 | 0.5 | 0.5 | 0 |
| | 0 | 0 | 0.5 | 1 |

| Referent | Property | | | | |
|----------|----------|-------|--------|------|-----|
| | square | green | circle | blue | N |
| | 0.94 | 0.06 | 0 | 0 | 144 |
| | 0 | 0.44 | 0.56 | 0 | 144 |
| | 0 | 0 | 0.17 | 0.83 | 144 |

Таблица 3: Результаты (Порождение)

3.2 Модель

$$EU_S(\text{wish to refer to } r, \text{ choose } p; \text{ parameter } f) = P_{\text{literal}} + f(p) \quad (30)$$

$$P_{\text{prod}}(\text{choose } p \mid \text{wish to refer to } r; \text{ parameters } \lambda, f) = \frac{e^{\lambda \cdot EU_S(r,p;f)}}{\sum_{p'} e^{\lambda \cdot EU_S(r,p';f)}} \quad (31)$$

$$P_{\text{comp}}(\text{choose } r \mid \text{receive } p; \text{ parameters } \lambda, f) = \frac{P(r) \cdot P_{\text{prod}}(p|r; \lambda, f)}{\sum_{r'} P(r') \cdot P_{\text{prod}}(p|r'; \lambda, f)} \quad (32)$$

$P(r)$ - степень «явности» объекта (определяется экспериментально)

3.2.1 Наилучшее совпадение

production: $\lambda = 4.13, x = -0.1$

comprehension: $\lambda = 3.46, x = -0.23$

$f(p) = x$ - для цвета

| Property | Referent | | | N |
|----------|----------|------|------|-----|
| | | | | |
| none | 0.3 | 0.12 | 0.58 | 240 |
| green | 0.36 | 0.64 | 0 | 180 |
| circle | 0 | 0.35 | 0.65 | 180 |

Таблица 4: Результаты (Понимание)

3.3 Скалярные импликатуры: RSA-модель

$$Pr_{R_0}(w|\varphi) = \begin{cases} \frac{1}{|\|\varphi\|_W|}, & \text{if } w \in \|\varphi\|_W \\ 0, & \text{else} \end{cases} \quad (33)$$

$$Pr_{S_1}(\varphi|w) = \frac{e^{\lambda \times \ln(Pr_{R_0}(w|\varphi))}}{\sum_{\varphi_i \in F} e^{\lambda \times \ln(Pr_{R_0}(w|\varphi_i))}} \quad (34)$$

$$Pr_{R_1}(w|\varphi) = \frac{Pr_{S_1}(\varphi|w) \times Pr(w)}{\sum_{w_i \in W} Pr_{S_1}(\varphi|w_i) \times Pr(w)} \quad (35)$$

Рис. 1: Pr_{R_0}

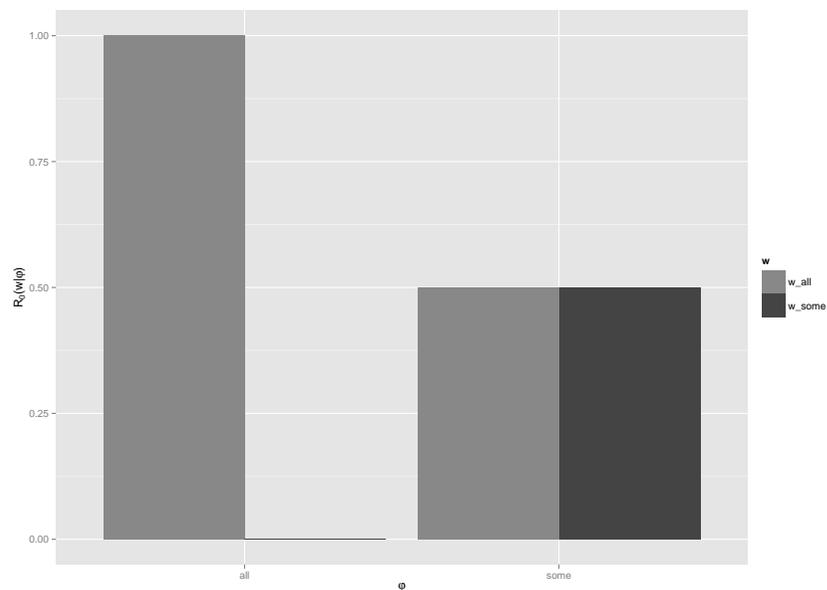


Рис. 2: $Pr_{R_1}, \lambda = 1$

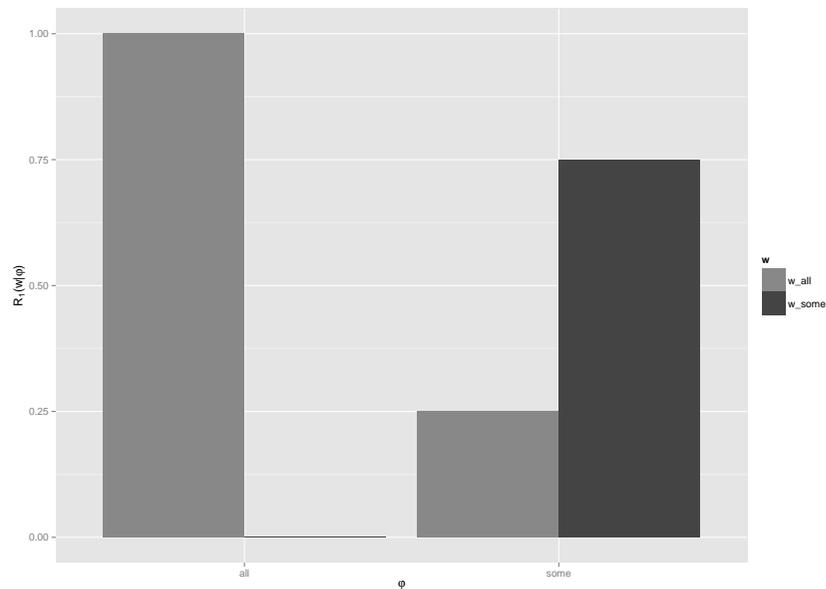


Рис. 3: $Pr_{R_1}, \lambda = 3$

